# Git & github

**Hadley Wickham**
@hadleywickham
Chief Scientist, RStudio

**February 2015**

# Overview

- Focus mostly on the why, not the how

- Learning git & github is definitely frustrating at first

- Two big pay offs: **safety** & **community**

- Mechanics at http://r-pkgs.had.co.nz/git.html

# Safety

"Failures, repeated failures, are finger posts on the road to achievement. One fails forward toward success."
— C. S. Lewis

data_2010.05.28-test
data_2010.05.28-retest
data_2010.05.28-re-retest
data_2010.05.28-calibrate
data_2010.05.29-aaaaargh
data_2010.05.29-##$#$!
data_2010.05.30-woohoo!
data_2010.05.30-USETHISONE

**Time**

The history of a project can be viewed as a series of changes

# Changes

- A unique identifier

- What changed?

- When did it change?

- Who changed it?

- Why did it change?

# Changes

- **A unique identifier**
- What changed?
- **When did it change?**
- Who changed it?
- **Why did it change?**

(And this obviously gets hard to coordinate for multiple files)

With git, each change (**commit**)
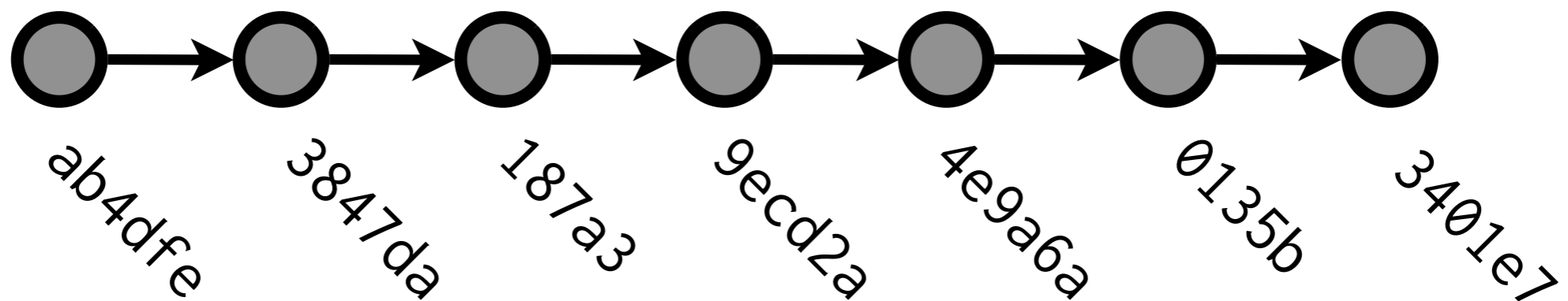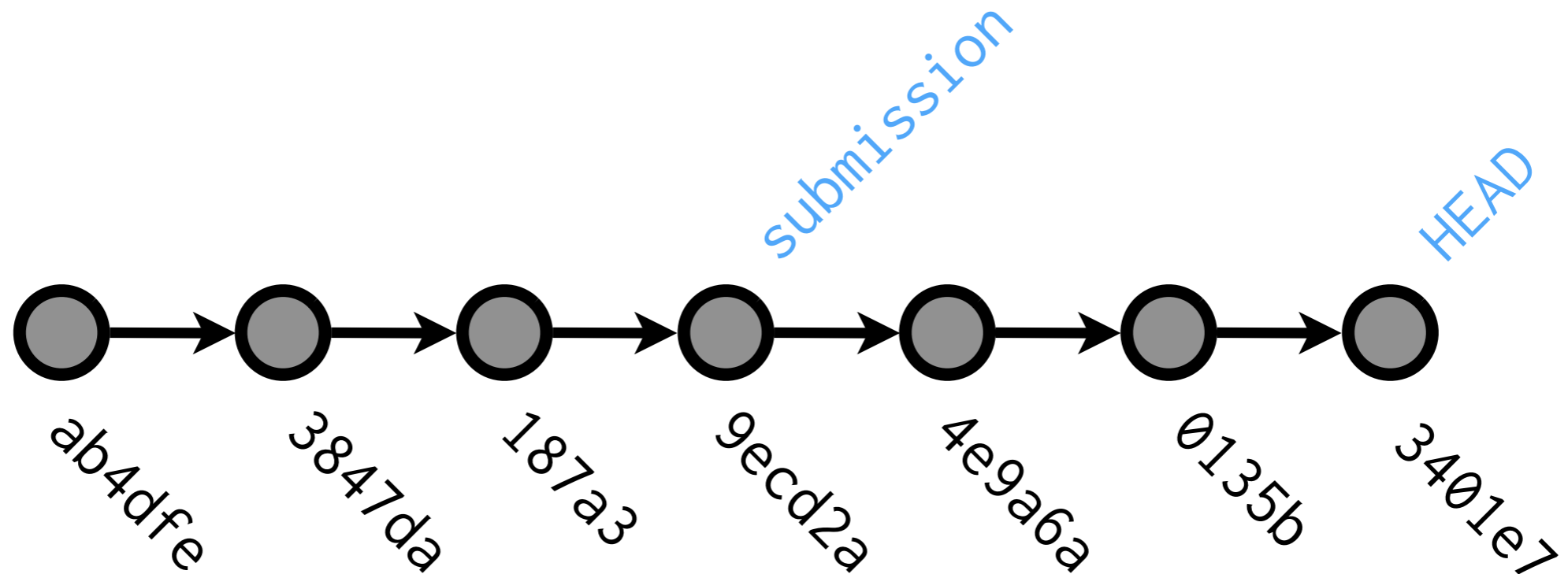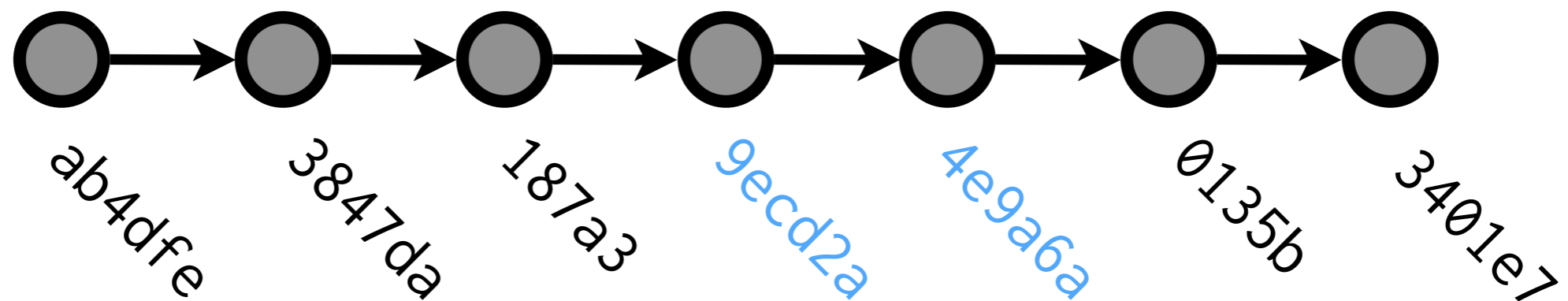is given a unique identifier, called
a **sha**



The sha is a key into a database
that provides the author, date,
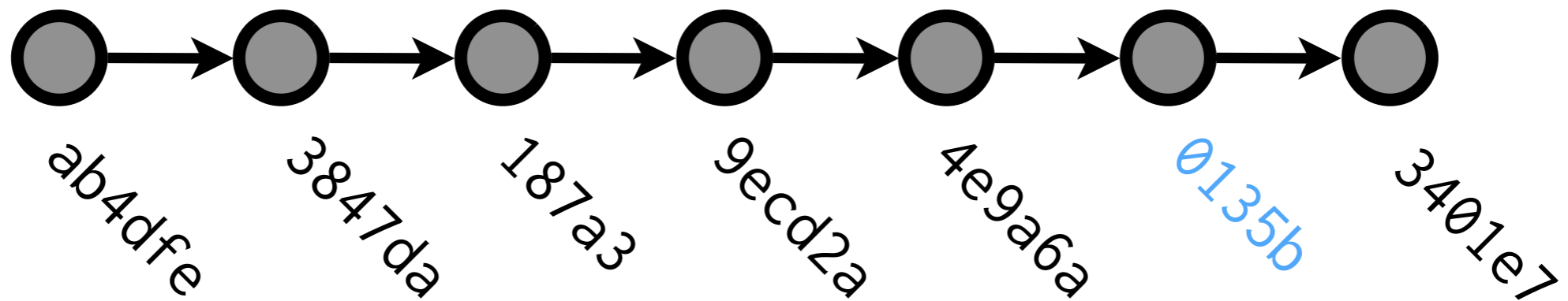and a description

# You can also name individual commits



```
git tag submission 9ecd2a
```
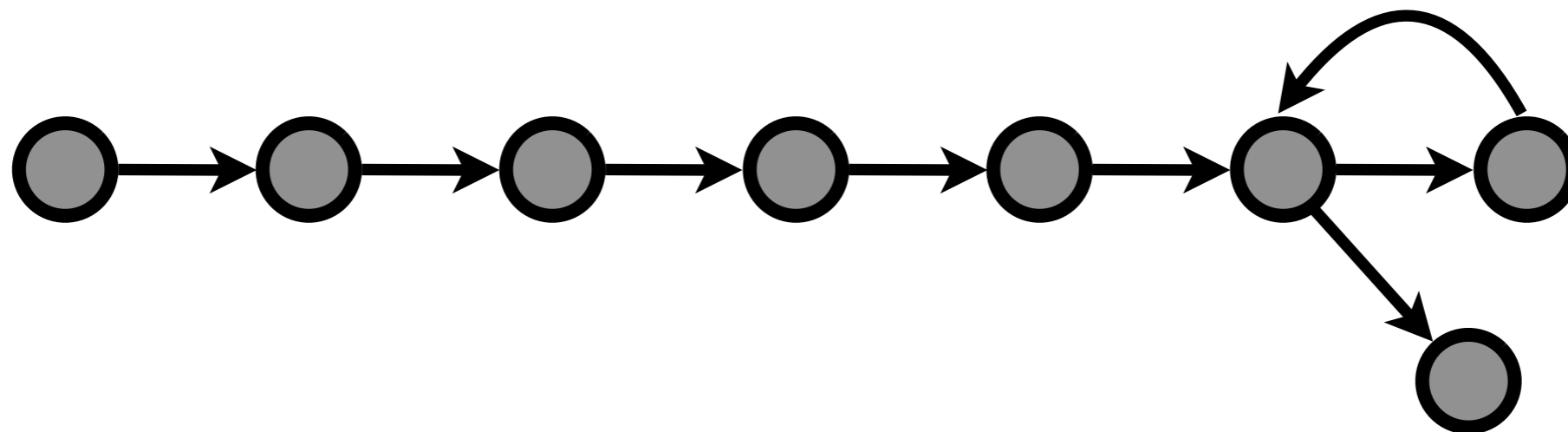
# Then see exactly what's changed



git diff 9ecd2a..4e9a6a
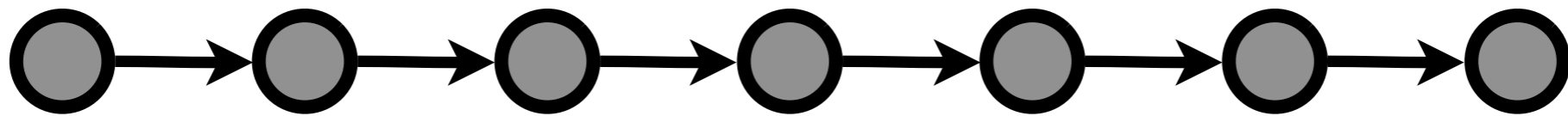
You can revert to a previous
change with git checkout
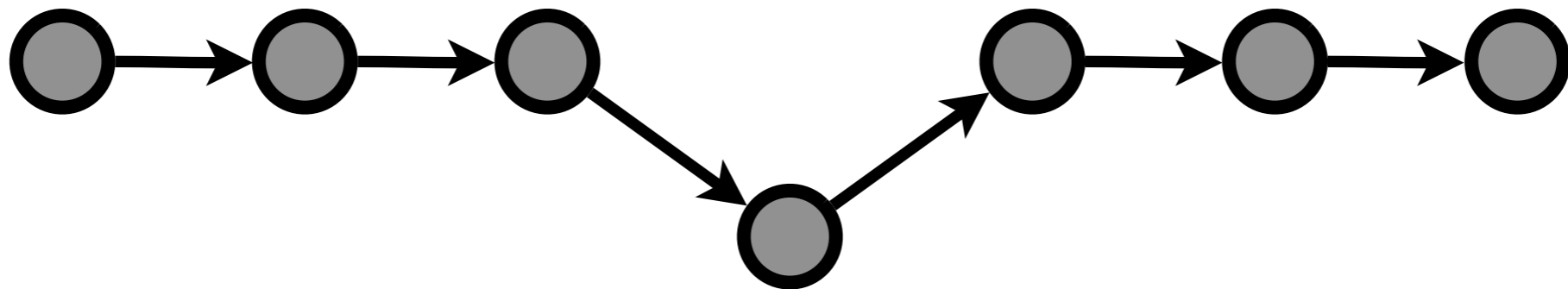

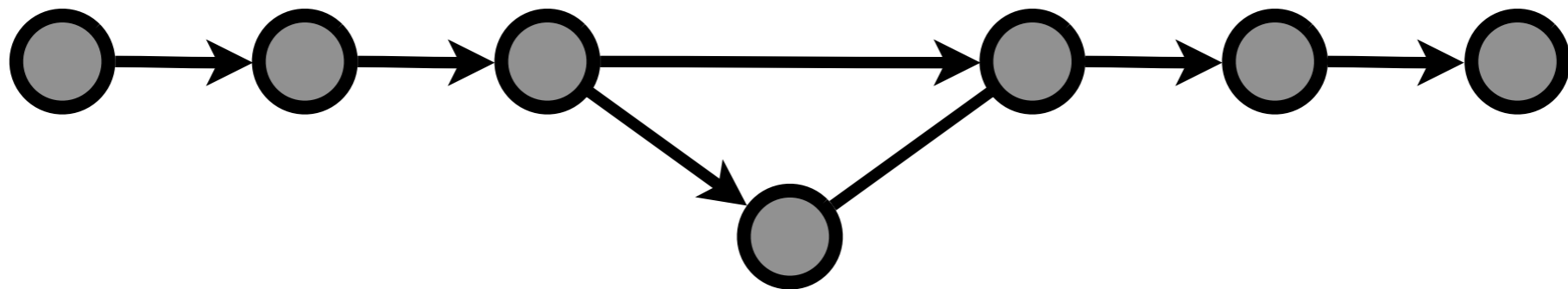
git checkout 0135b

That allows you to undo
mistakes

Alternatively you can go back in time and pretend mistakes never existed

Alternatively you can go back in time and pretend mistakes never existed

Alternatively you can go back in time and pretend mistakes never existed. **That's risky!**



git rebase ...

# Demo

http://r-pkgs.had.co.nz
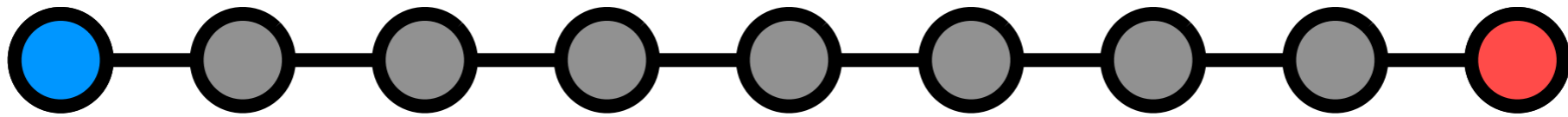
http://github.com/hadley/r-pkgs

# Demo

Vows package & roxygen2

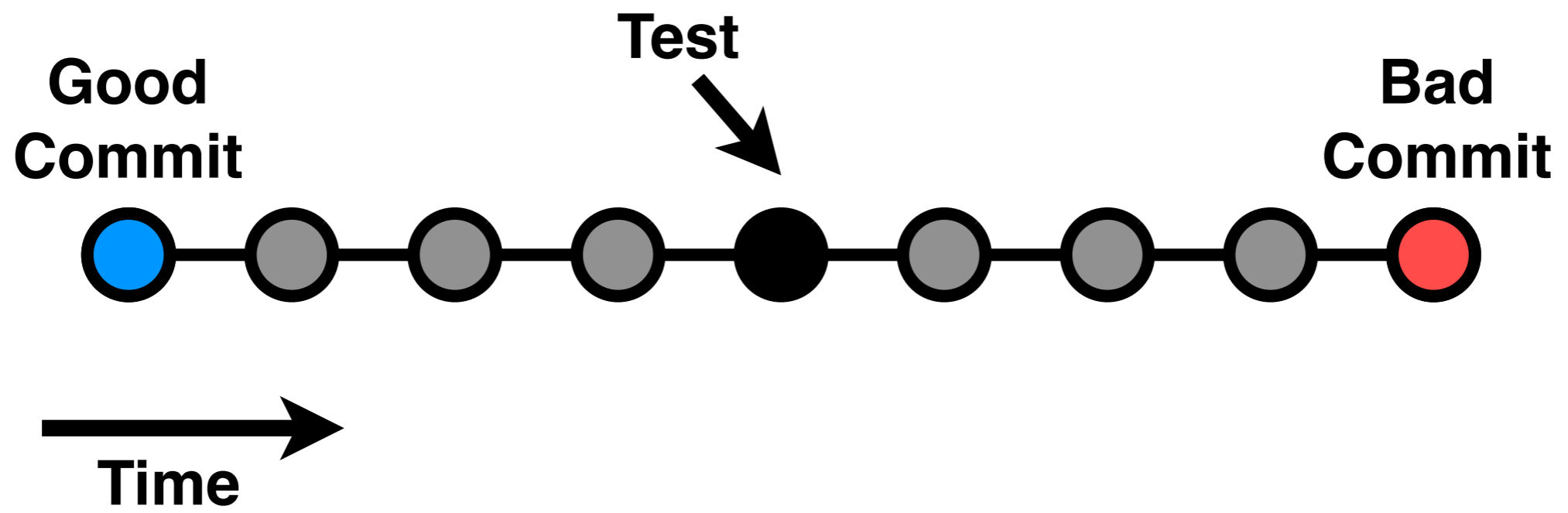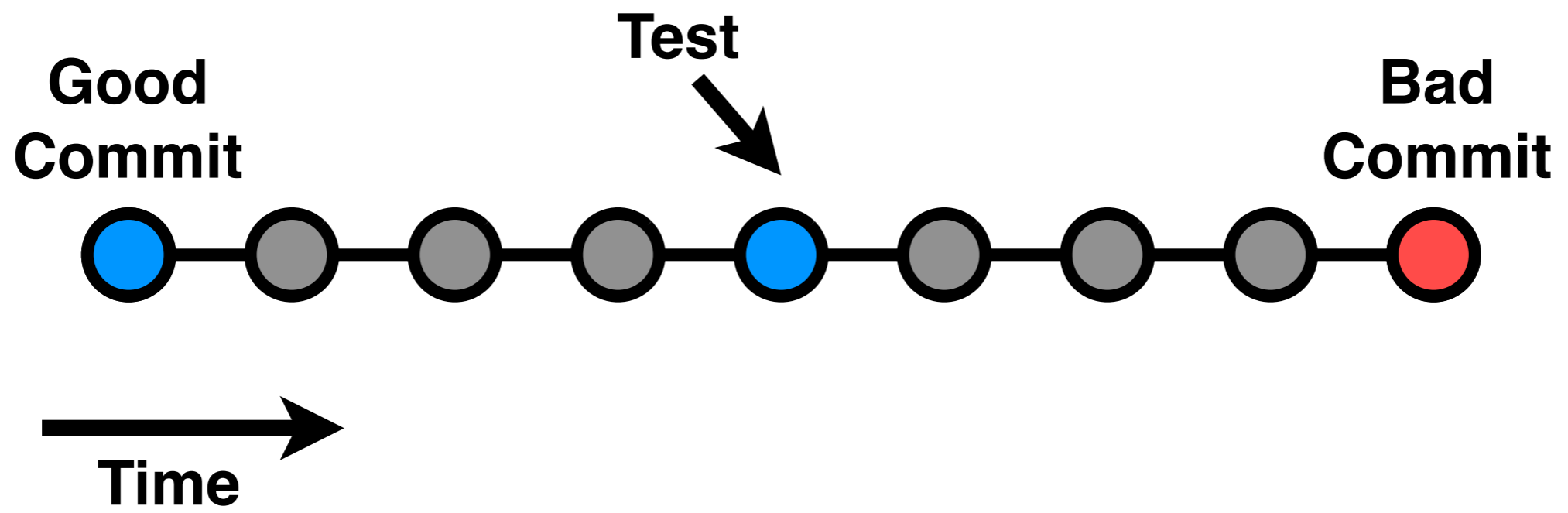# Other benefits

git bisect

https://github.com/wch/bisectr

**Good Commit**

**Bad Commit**

**Time**
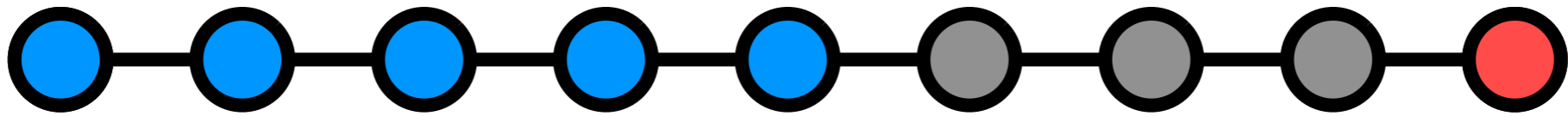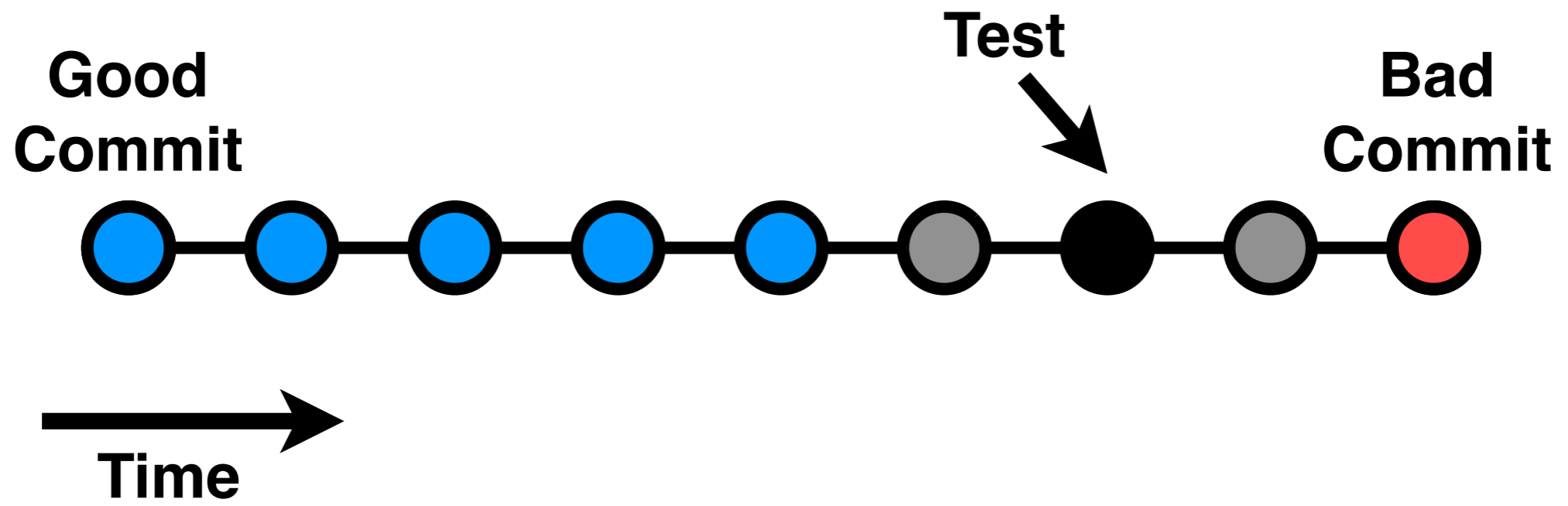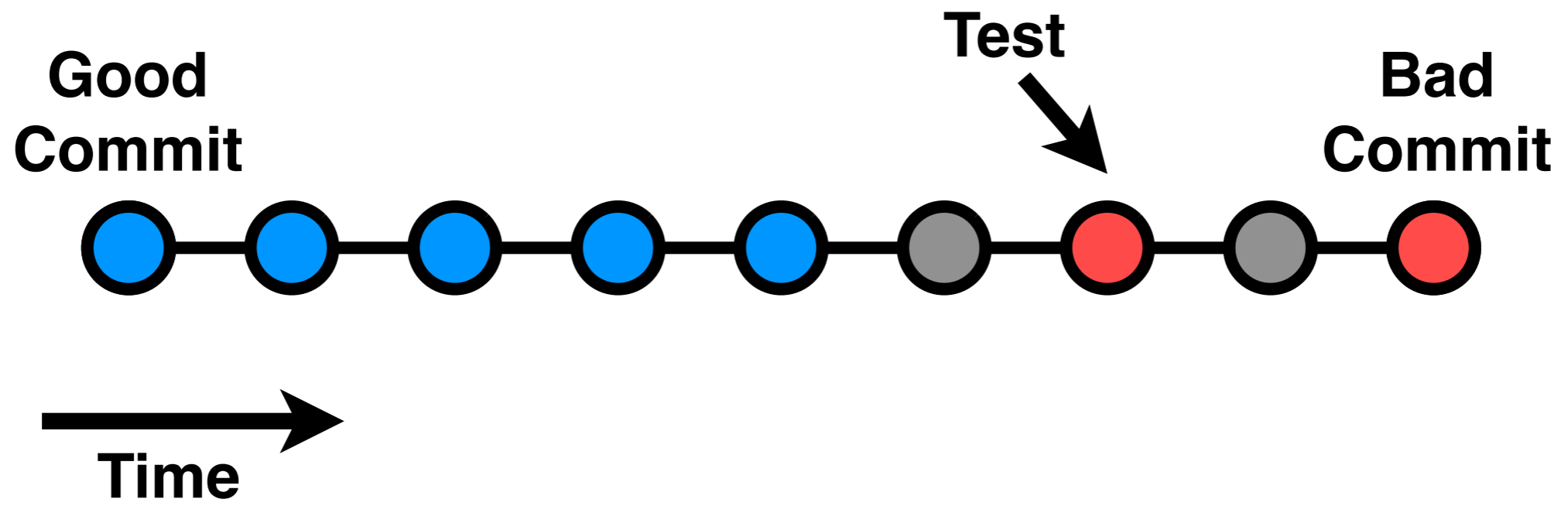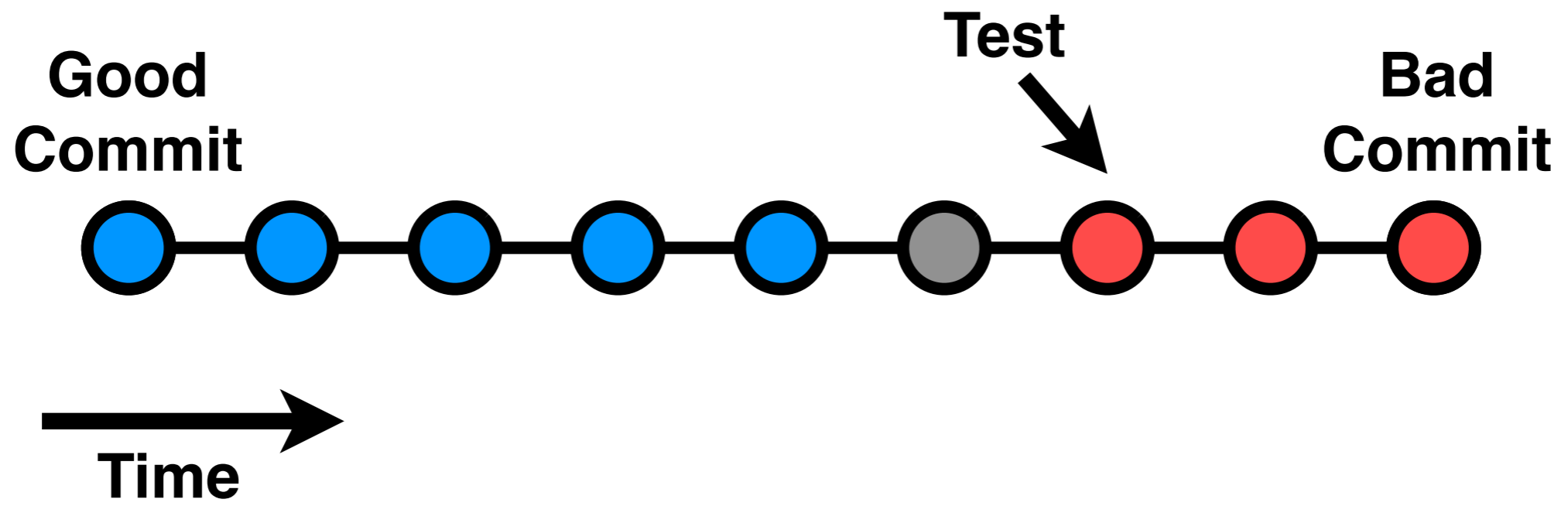
**Where did the problem happen?**

Good Commit

Test

Bad Commit

Time

# Community

# A home on the web

# Nicely formatted readme



**README.md**

# dplyr

build failing

dplyr is the next iteration of plyr, focussed on tools for working with data frames (hence the `d` in the name). It has three main goals:

- Identify the most important data manipulation tools needed for data analysis and make them easy to use from R.

- Provide blazing fast performance for in-memory data by writing key pieces in C++.

- Use the same interface to work with data no matter where it's stored, whether in a data frame, a data table or database.

`dplyr` is not yet available on CRAN, but you can install it from github with:

```
devtools::install_github("dplyr")
```

To get started, read the notes below, then read the intro vignette:
`vignette("introduction", package = "dplyr")`. To make the most of dplyr, I also recommend that you familiarise yourself with the principles of tidy data: this will help you get your data into a form that works well with dplyr, ggplot2 and R's many modelling functions.

If you encounter a clear bug, please file a minimal reproducible example on github. For questions and other discussion, please use the manipulatr mailing list.

```
# Easy to install a package from github
library(devtools)
install_github("hadley/devtools")

# Github is a viable location for package
# distribution, but it doesn't yet have the
# reach or authority of CRAN. Currently
# about ~3,000 packages on github.
```

# See what's happening

# See exactly what changed

# Track issues

# Talk things over

GitHub, Inc. [US] https://github.com/hadley/dplyr/issues/123

**hadley** opened this issue 21 days ago    Edit

## C++ code needs to regularly check for user interupts

**Open**

4 comments

**romainfrancois** is assigned ⚙▾         Milestone: **v0.1**    ⚙▾

Labels  ⚙▾

So if you accidentally miss-specify a call you can abort it.

enhancement

2 participants

◄ **romainfrancois** referenced this issue from a commit                14 days ago

**romainfran...** adding code to check for interupts. #123        ✕ 💬 0627028

💬 **romainfrancois** commented                                      14 days ago  ✏ ⊗

I'm adding calls to `R_CheckUserInterrupt` here and there. But I'm wondering if this plays well with C++.

I think this uses long jumps and that possibly bypasses C++ destructors. I might ask this on R-devel.

💬 **romainfrancois** commented                                      14 days ago  ✏ ⊗

This comes from Simon's contribution to a thread on R-devel :

```
static void chkIntFn(void *dummy) {
  R_CheckUserInterrupt();
}

// this will call the above in a top-level context so it won't longjmp-out of your context
bool checkInterrupt() {
  return (R_ToplevelExec(chkIntFn, NULL) == FALSE);
}

// your code somewhere ...
```

# Receive pull requests

# … and review online

# Find out as soon as you break something

# Beautiful release notes

# Demo

Easily contribute small fixes

https://github.com/mjockers/syuzhet

# My advice

# Learning Git & Github

- Is frustrating at first & there's a lot of new vocab to learn.

- Pays off by increasing the pace of your coding because you don't have to worry as much about mistakes.

- Pays off by connecting you with a community of like minded people who can help make your code better.

http://r-pkgs.had.co.nz/git.html